# Application of Information Theory, Lecture 7
## Relative Entropy

Iftach Haitner

Tel Aviv University.

December 1, 2015

# Part I

# **Statistical Distance**

# Statistical distance

## Statistical distance

- Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

## Statistical distance

- Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$
- Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

# Statistical distance

- Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

- Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

- This is simply the $L_1$ norm between the distribution vectors

## Statistical distance

- ▶ Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$
- ▶ Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

- ▶ This is simply the $L_1$ norm between the distribution vectors
- ▶ We will see other "distance" measures for distributions next lecture

## Statistical distance

▸ Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

▸ Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

▸ This is simply the $L_1$ norm between the distribution vectors

▸ We will see other "distance" measures for distributions next lecture

▸ For $Z \sim p$ and $Y \sim q$, let $\mathrm{SD}(X, Y) = \mathrm{SD}(p, q)$

# Statistical distance

- ▶ Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

- ▶ Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

- ▶ This is simply the $L_1$ norm between the distribution vectors

- ▶ We will see other "distance" measures for distributions next lecture

- ▶ For $Z \sim p$ and $Y \sim q$, let $\mathrm{SD}(X, Y) = \mathrm{SD}(p, q)$

- ▶ Claim (HW): $\mathrm{SD}(p, q) = \max_{\mathcal{S} \subseteq [m]} \left( \sum_{i \in \mathcal{S}} p_i - \sum_{i \in \mathcal{S}} q_i \right)$

# Statistical distance

- Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

- Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

- This is simply the $L_1$ norm between the distribution vectors

- We will see other "distance" measures for distributions next lecture

- For $Z \sim p$ and $Y \sim q$, let $\mathrm{SD}(X, Y) = \mathrm{SD}(p, q)$

- Claim (HW): $\mathrm{SD}(p, q) = \max_{S \subseteq [m]} \left( \sum_{i \in S} p_i - \sum_{i \in S} q_i \right)$

- Hence, $\mathrm{SD}(p, q) = \max_{\mathrm{D}} \left( \Pr_{X \sim p}[\mathrm{D}(X) = 1] - \Pr_{X \sim q}[\mathrm{D}(X) = 1] \right)$

# Statistical distance

- Let $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$ be distributions over $[m]$

- Their statistical distance (also known as, variation distance) is defined by

$$\mathrm{SD}(p, q) := \frac{1}{2} \sum_{i \in [m]} |p_i - q_i|$$

- This is simply the $L_1$ norm between the distribution vectors

- We will see other "distance" measures for distributions next lecture

- For $Z \sim p$ and $Y \sim q$, let $\mathrm{SD}(X, Y) = \mathrm{SD}(p, q)$

- Claim (HW): $\mathrm{SD}(p, q) = \max_{\mathcal{S} \subseteq [m]} \left( \sum_{i \in \mathcal{S}} p_i - \sum_{i \in \mathcal{S}} q_i \right)$

- Hence, $\mathrm{SD}(p, q) = \max_{\mathrm{D}} \left( \Pr_{X \sim p} [\mathrm{D}(X) = 1] - \Pr_{X \sim q} [\mathrm{D}(X) = 1] \right)$

- Interpretation

**Distance from the uniform distribution**

# Distance from the uniform distribution

- Let $X$ be rv over $[m]$

## Distance from the uniform distribution

- Let $X$ be rv over $[m]$
- $H(X) \leq \log m$

# Distance from the uniform distribution

- Let $X$ be rv over $[m]$
- $H(X) \leq \log m$
- $H(X) = \log m \longleftrightarrow X$ is uniform over $[m]$

# Distance from the uniform distribution

- Let $X$ be rv over $[m]$
- $H(X) \leq \log m$
- $H(X) = \log m \longleftrightarrow X$ is uniform over $[m]$

**Theorem 1 (this lecture)**

*Let $X$ rv over $[m]$. Assume $H(X) \geq \log m - \varepsilon$, then*

$$\mathrm{SD}(X, \sim [m]) \leq \sqrt{\varepsilon \cdot \frac{\ln 2}{2}} = O(\sqrt{\varepsilon})$$

# Part II

# **Relative entropy Distance**

Section 1

**Definition and Basic Facts**

## Definition

- For $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$, let

$$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$$

$0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$

## Definition

▶ For $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$, let

$$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$$

$0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$

▶ The relative entropy of pair of rv's, is the relative entropy of their distributions.

## Definition

▶ For $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$, let

$$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$$

$0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$

▶ The relative entropy of pair of rv's, is the relative entropy of their distributions.

▶ Names: Entropy of $p$ relative to $q$, relative entropy, information divergence, Kullback-Leibler (KL) divergence/distance

## Definition

- For $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$, let

$$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$$

$0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$

- The relative entropy of pair of rv's, is the relative entropy of their distributions.

- Names: Entropy of $p$ relative to $q$, relative entropy, information divergence, Kullback-Leibler (KL) divergence/distance

- Many different interpretations

## Definition

▶ For $p = (p_1, \ldots, p_m)$ and $q = (q_1, \ldots, q_m)$, let

$$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$$

$0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$

▶ The relative entropy of pair of rv's, is the relative entropy of their distributions.

▶ Names: Entropy of $p$ relative to $q$, relative entropy, information divergence, Kullback-Leibler (KL) divergence/distance

▶ Many different interpretations

▶ Main interpretation: the information we gained about $X$, if we originally thought $X \sim q$ and now we learned $X \sim p$

**Numerical Example**

# Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$
- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$
- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$

- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0 = \frac{1}{4} \cdot (-1) + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1 = \frac{1}{2}$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$

- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0 = \frac{1}{4} \cdot (-1) + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1 = \frac{1}{2}$

- $D(q\|p) = \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{8} \log \frac{\frac{1}{8}}{\frac{1}{4}} + \frac{1}{8} \log \frac{\frac{1}{8}}{0}$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$

- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0 = \frac{1}{4} \cdot (-1) + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1 = \frac{1}{2}$

- $D(q\|p) = \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{8} \log \frac{\frac{1}{8}}{\frac{1}{4}} + \frac{1}{8} \log \frac{\frac{1}{8}}{0}$

## Numerical Example

$D(p\|q) = \sum_{i=1}^{m} p_i \log \frac{p_i}{q_i}$

- $p = (\frac{1}{4}, \frac{1}{2}, \frac{1}{4}, 0)$, $q = (\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8})$

- $D(p\|q) = \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{8}} + 0 \log 0 = \frac{1}{4} \cdot (-1) + \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1 = \frac{1}{2}$

- $D(q\|p) = \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{4}} + \frac{1}{4} \log \frac{\frac{1}{4}}{\frac{1}{2}} + \frac{1}{8} \log \frac{\frac{1}{8}}{\frac{1}{4}} + \frac{1}{8} \log \frac{\frac{1}{8}}{0} = $
  $\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot (-1) + \frac{1}{8} \cdot (-1) + \infty = \infty$

# Supporting the interpretation

## Supporting the interpretation

- $X$ rv over $[m]$

## Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$

## Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$
  (just by knowing its distribution)

# Supporting the interpretation

- ▶ $X$ rv over $[m]$
- ▶ $H(X)$ — measure for amount of information we do not have about $X$
- ▶ $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- ▶ Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$

## Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$

## Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
- 
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$
  (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
- 
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$

## Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$
  (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
-
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$
$$= \log m + \sum_i p_i \log p_i$$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$
  (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
- 
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$
$$= \log m + \sum_i p_i \log p_i = \sum_i p_i (\log p_i - \log \frac{1}{m})$$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
- 
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$
$$= \log m + \sum_i p_i \log p_i = \sum_i p_i (\log p_i - \log \frac{1}{m})$$
$$= \sum_i p_i \log \frac{p_i}{\frac{1}{m}}$$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$ (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
-
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$
$$= \log m + \sum_i p_i \log p_i = \sum_i p_i (\log p_i - \log \frac{1}{m})$$
$$= \sum_i p_i \log \frac{p_i}{\frac{1}{m}} = D(p \| \sim [m])$$

# Supporting the interpretation

- $X$ rv over $[m]$
- $H(X)$ — measure for amount of information we do not have about $X$
- $\log m - H(X)$ — measure for information we do have about $X$
  (just by knowing its distribution)
- Example $X = (X_1, X_2) \sim (\frac{1}{2}, 0, 0, \frac{1}{2})$ over $\{00, 01, 10, 11\}$
- $H(X) = 1$, $\log m - H(X) = 2 - 1 = 1$
- Indeed, we know $X_1 \oplus X_2$
-
$$H(\sim [m]) - H(p_1, \ldots, p_m) = \log m - H(p_1, \ldots, p_m)$$
$$= \log m + \sum_i p_i \log p_i = \sum_i p_i (\log p_i - \log \frac{1}{m})$$
$$= \sum_i p_i \log \frac{p_i}{\frac{1}{m}} = D(p \| \sim [m])$$
- $D(X \| \sim [m])$ — measures the information we gained about $X$, if we
  originally thought it is $\sim [m]$ and now we learned it is $\sim p$

# Supporting the interpretation, cont.

## Supporting the interpretation, cont.

- (generally) $D(p\|q) \neq H(p) - H(q)$

# Supporting the interpretation, cont.

- (generally) $D(p\|q) \neq H(p) - H(q)$
- $H(q) - H(p)$ is not a good measure for information change

# Supporting the interpretation, cont.

- (generally) $D(p\|q) \neq H(p) - H(q)$
- $H(q) - H(p)$ is not a good measure for information change
- Example: $q = (0.01, 0.99)$ and $p = (0.99, 0.01)$

# Supporting the interpretation, cont.

- ▶ (generally) $D(p\|q) \neq H(p) - H(q)$
- ▶ $H(q) - H(p)$ is not a good measure for information change
- ▶ Example: $q = (0.01, 0.99)$ and $p = (0.99, 0.01)$
- ▶ We were almost sure that $X = 1$ but learned that $X$ is almost surely 0

# Supporting the interpretation, cont.

- (generally) $D(p\|q) \neq H(p) - H(q)$

- $H(q) - H(p)$ is not a good measure for information change

- Example: $q = (0.01, 0.99)$ and $p = (0.99, 0.01)$

- We were almost sure that $X = 1$ but learned that $X$ is almost surely 0

- But $H(p) - H(q) = 0$

# Supporting the interpretation, cont.

- ▶ (generally) $D(p\|q) \neq H(p) - H(q)$

- ▶ $H(q) - H(p)$ is not a good measure for information change

- ▶ Example: $q = (0.01, 0.99)$ and $p = (0.99, 0.01)$

- ▶ We were almost sure that $X = 1$ but learned that $X$ is almost surely 0

- ▶ But $H(p) - H(q) = 0$

- ▶ Also, $H(q) - H(p)$ might be negative

# Supporting the interpretation, cont.

- (generally) $D(p\|q) \neq H(p) - H(q)$

- $H(q) - H(p)$ is not a good measure for information change

- Example: $q = (0.01, 0.99)$ and $p = (0.99, 0.01)$

- We were almost sure that $X = 1$ but learned that $X$ is almost surely 0

- But $H(p) - H(q) = 0$

- Also, $H(q) - H(p)$ might be negative

- We understand $D(p\|q)$ as the information we gained about $X$, if we originally thought it is $\sim q$ and now we learned it is $\sim p$

## Changing distribution

► What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?
How can a distribution change?

## Changing distribution

- ▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

  How can a distribution change?

- ▶ Typically, this happens by learning additional information

## Changing distribution

▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

How can a distribution change?

▶ Typically, this happens by learning additional information

▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i | E]$

## Changing distribution

- What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

  How can a distribution change?

- Typically, this happens by learning additional information

- $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i|E]$

- Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

# Changing distribution

- ▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

  How can a distribution change?

- ▶ Typically, this happens by learning additional information

- ▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i|E]$

- ▶ Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

- ▶ The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

# Changing distribution

- ▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

  How can a distribution change?

- ▶ Typically, this happens by learning additional information

- ▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i | E]$

- ▶ Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

- ▶ The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

- ▶ Another example

| $X^Y$ | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | 0 | 0 |
| 1 | $\frac{1}{4}$ | 0 | $\frac{1}{4}$ | 0 |

## Changing distribution

▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

   How can a distribution change?

▶ Typically, this happens by learning additional information

▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i | E]$

▶ Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

▶ The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

▶ Another example

| X\Y | 1 | 2 | 3 | 4 |
|-----|---|---|---|---|
| 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | 0 | 0 |
| 1 | $\frac{1}{4}$ | 0 | $\frac{1}{4}$ | 0 |

▶ $Y \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$, but

# Changing distribution

▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

   How can a distribution change?

▶ Typically, this happens by learning additional information

▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i | E]$

▶ Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

▶ The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

▶ Another example

| X\Y | 1 | 2 | 3 | 4 |
|-----|---|---|---|---|
| 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | 0 | 0 |
| 1 | $\frac{1}{4}$ | 0 | $\frac{1}{4}$ | 0 |

▶ $Y \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$, but

▶ $Y \sim (\frac{1}{2}, \frac{1}{2}, 0, 0)$ conditioned on $X = 0$

# Changing distribution

▶ What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

   How can a distribution change?

▶ Typically, this happens by learning additional information

▶ $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i|E]$

▶ Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

▶ The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

▶ Another example

| X\Y | 1 | 2 | 3 | 4 |
|-----|---|---|---|---|
| 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | 0 | 0 |
| 1 | $\frac{1}{4}$ | 0 | $\frac{1}{4}$ | 0 |

▶ $Y \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$, but

▶ $Y \sim (\frac{1}{2}, \frac{1}{2}, 0, 0)$ conditioned on $X = 0$

▶ $Y \sim (\frac{1}{2}, 0, \frac{1}{2}, 0)$ conditioned on $X = 1$

# Changing distribution

► What does it mean: originally thought $X \sim q$ and now we learned $X \sim p$?

   How can a distribution change?

► Typically, this happens by learning additional information

► $q_i = \Pr[X = i]$ and $p_i = \Pr[X = i | E]$

► Example $X \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$; someone saw $X$ and tells us that $X \leq 2$

► The distribution changes to $X \sim (\frac{2}{3}, \frac{1}{3}, 0, 0)$

► Another example

| $X^Y$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0 | $\frac{1}{4}$ | $\frac{1}{4}$ | 0 | 0 |
| 1 | $\frac{1}{4}$ | 0 | $\frac{1}{4}$ | 0 |

► $Y \sim (\frac{1}{2}, \frac{1}{4}, \frac{1}{4}, 0)$, but

► $Y \sim (\frac{1}{2}, \frac{1}{2}, 0, 0)$ conditioned on $X = 0$

► $Y \sim (\frac{1}{2}, 0, \frac{1}{2}, 0)$ conditioned on $X = 1$

► Generally, a distribution can change if we condition on event $E$

# Additional properties

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$
- $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p\|q) = \infty$

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$
- $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p \| q) = \infty$
- If originally $\Pr[X = i] = 0$, then it cannot be more than 0 after we learned something.

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$

- $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p \| q) = \infty$

- If originally $\Pr[X = i] = 0$, then it cannot be more than $0$ after we learned something.

- Hence, it make sense to think of it as infinite amount of information learnt

## Additional properties

- ▶ $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$

- ▶ $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p\|q) = \infty$

- ▶ If originally $\Pr[X = i] = 0$, then it cannot be more than $0$ after we learned something.

- ▶ Hence, it make sense to think of it as infinite amount of information learnt

- ▶ Alteratively, we can define $D(p\|q)$ only for distribution with
  $q_i = 0 \implies p_i = 0$
  (recall that $\Pr[X = i] = 0 \implies \Pr[X = i|E] = 0$, for any event $E$

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$

- $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p\|q) = \infty$

- If originally $\Pr[X = i] = 0$, then it cannot be more than $0$ after we learned something.

- Hence, it make sense to think of it as infinite amount of information learnt

- Alteratively, we can define $D(p\|q)$ only for distribution with
  $q_i = 0 \implies p_i = 0$
  (recall that $\Pr[X = i] = 0 \implies \Pr[X = i|E] = 0$, for any event $E$

- If $p_i$ is large and $q_i$ is small, then $D(p\|q)$ is large

## Additional properties

- $0 \log \frac{0}{0} = 0$, $p \log \frac{p}{0} = \infty$ for $p > 0$

- $\exists i$ s.t. $p_i > 0$ and $q_i = 0$, then $D(p\|q) = \infty$

- If originally $\Pr[X = i] = 0$, then it cannot be more than $0$ after we learned something.

- Hence, it make sense to think of it as infinite amount of information learnt

- Alteratively, we can define $D(p\|q)$ only for distribution with
  $q_i = 0 \implies p_i = 0$
  (recall that $\Pr[X = i] = 0 \implies \Pr[X = i|E] = 0$, for any event $E$

- If $p_i$ is large and $q_i$ is small, then $D(p\|q)$ is large

- $D(p\|q) \geq 0$, with equality iff $p = q$ (hw)

# Example

## Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^{n} q_i = 2^{-k}$ (i.e., $n < m$)

## Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^{n} q_i = 2^{-k}$ (i.e., $n < m$)
- $p_i = \begin{cases} q_i/2^{-k}, & 1 \leq i \leq n \\ 0, & \text{otherwise.} \end{cases}$

# Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^{n} q_i = 2^{-k}$ (i.e., $n < m$)

- $p_i = \begin{cases} q_i/2^{-k}, & 1 \leq i \leq n \\ 0, & \text{otherwise.} \end{cases}$

- $p = (p_1, \ldots, p_m)$ — the distribution of $q$ conditioned on the event $i \in [n]$

## Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^{n} q_i = 2^{-k}$ (i.e., $n < m$)

- $p_i = \begin{cases} q_i/2^{-k}, & 1 \leq i \leq n \\ 0, & \text{otherwise.} \end{cases}$

- $p = (p_1, \ldots, p_m)$ — the distribution of $q$ conditioned on the event $i \in [n]$

- $D(p\|q) = \sum_{i=1}^{n} p_i \log \frac{p_i}{q_i} = \sum_{i=1}^{n} p_i \log 2^k = \sum_{i=1}^{n} p_i k = k$

## Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^{n} q_i = 2^{-k}$ (i.e., $n < m$)
- $p_i = \begin{cases} q_i/2^{-k}, & 1 \leq i \leq n \\ 0, & \text{otherwise.} \end{cases}$
- $p = (p_1, \ldots, p_m)$ — the distribution of $q$ conditioned on the event $i \in [n]$
- $D(p \| q) = \sum_{i=1}^{n} p_i \log \frac{p_i}{q_i} = \sum_{i=1}^{n} p_i \log 2^k = \sum_{i=1}^{n} p_i k = k$
- We gained $k$ bits of information

# Example

- $q = (q_1, \ldots, q_m)$ with $\sum_{i=1}^n q_i = 2^{-k}$ (i.e., $n < m$)

- $p_i = \begin{cases} q_i/2^{-k}, & 1 \leq i \leq n \\ 0, & \text{otherwise.} \end{cases}$

- $p = (p_1, \ldots, p_m)$ — the distribution of $q$ conditioned on the event $i \in [n]$

- $D(p\|q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i} = \sum_{i=1}^n p_i \log 2^k = \sum_{i=1}^n p_i k = k$

- We gained $k$ bits of information

- Example: $\sum_{i=1}^n q_i = \frac{1}{2}$, and we were told that $i \leq n$ or $i > n$, we got one bit of information

Section 2

**Axiomatic Derivation**

# Axiomatic derivation

# Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$

2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1 - \alpha) p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1 - \alpha) q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p\| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m)\|(q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m)\|(q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

# Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) = \tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1 - \alpha) p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1 - \alpha) q_m))$, for any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

▶ $\tilde{D}(p \| q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m) \|$
   $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p \| q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m) \|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$,

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$

2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p \| q) = D((\alpha_{1,1} p_1, \ldots, \alpha_{1,k_1} p_1, \ldots, \alpha_{m,1} p_m, \ldots, \alpha_{m,k_m} p_m) \|$
  $(\alpha_{1,1} q_1, \ldots, \alpha_{1,k_1} q_1, \ldots, \alpha_{m,1} q_m, \ldots, \alpha_{m,k_m} q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$,

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p \| \sim [m]) = \log m - H(p)$

2. $\tilde{D}((p_1, \ldots, p_m) \| (q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m) \| (q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p \| q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m) \|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$, it follows that

$\tilde{D}(p \| q) = \log M - H((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m))$

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p\| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m)\|(q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m)\|(q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p\|q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m)\|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$, it follows that

$\tilde{D}(p\|q) = \log M - H((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m))$

$= \sum p_i \log M + \sum_i p_i \log \alpha_i p_i$

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p\| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m)\|(q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1 - \alpha)p_m)\|(q_1, \ldots, q_{m-1}, \alpha q_m, (1 - \alpha)q_m))$, for any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p\|q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m)\|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$, it follows that

$$\tilde{D}(p\|q) = \log M - H((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m))$$
$$= \sum p_i \log M + \sum_i p_i \log \alpha_i p_i = \sum_i p_i (\log M + \log \frac{p_i}{q_i M})$$

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p\| \sim [m]) = \log m - H(p)$
2. $\tilde{D}((p_1, \ldots, p_m)\|(q_1, \ldots, q_m)) =$
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1-\alpha)p_m)\|(q_1, \ldots, q_{m-1}, \alpha q_m, (1-\alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p\|q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m)\|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j \geq 0}$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$, it follows that

$\tilde{D}(p\|q) = \log M - H((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m))$

$$= \sum p_i \log M + \sum_i p_i \log \alpha_i p_i = \sum_i p_i(\log M + \log \frac{p_i}{q_i M}) = \sum_i p_i \log \frac{p_i}{q_i}.$$

## Axiomatic derivation

Let $\tilde{D}$ is a continuous and symmetric (wrt each distribution) function such that

1. $\tilde{D}(p\| \sim [m]) = \log m - H(p)$

2. $\tilde{D}((p_1, \ldots, p_m)\|(q_1, \ldots, q_m)) = $
   $\tilde{D}((p_1, \ldots, p_{m-1}, \alpha p_m, (1 - \alpha)p_m)\|(q_1, \ldots, q_{m-1}, \alpha q_m, (1 - \alpha)q_m))$, for
   any $\alpha \in [0, 1]$

then $\tilde{D} = D$.

Interpretation

Proof:

- $\tilde{D}(p\|q) = D((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m)\|$
  $(\alpha_{1,1}q_1, \ldots, \alpha_{1,k_1}q_1, \ldots, \alpha_{m,1}q_m, \ldots, \alpha_{m,k_m}q_m))$, for $\sum_j \alpha_{i,j} = 1$ and $\alpha_{i,j} \geq 0$

- Taking $\alpha$'s s.t. $\alpha_{i,1} = \alpha_{i,2} \ldots, \alpha_{i,k_i} = \alpha_i$ and $\alpha_i q_i = \frac{1}{M}$, it follows that

$\tilde{D}(p\|q) = \log M - H((\alpha_{1,1}p_1, \ldots, \alpha_{1,k_1}p_1, \ldots, \alpha_{m,1}p_m, \ldots, \alpha_{m,k_m}p_m))$

$= \sum p_i \log M + \sum_i p_i \log \alpha_i p_i = \sum_i p_i(\log M + \log \frac{p_i}{q_i M}) = \sum_i p_i \log \frac{p_i}{q_i}.$

- Zeros and non-rational $q_i$'s are dealt by continuity

Section 3

**Relation to Mutual Information**

# Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$

# Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\qquad p_{0,i} = \Pr[X = i | Y = 0]$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$

- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$,      $p_{0,i} = \Pr[X = i|Y = 0]$

- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$,      $p_{1,i} = \Pr[X = i|Y = 1]$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\quad p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$, $\quad p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

# Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\qquad p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$, $\qquad p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$\underset{Y}{\mathrm{E}}[D(p_Y \| q)] = \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m)$$
$$+ \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m)$$

# Mutual information as expected relative entropy

- ▶ Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- ▶ $(X | Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$,      $p_{0,i} = \Pr[X = i | Y = 0]$
- ▶ $(X | Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$,      $p_{1,i} = \Pr[X = i | Y = 1]$
- ▶ If we learned $Y = j$, we gained $D(p_j \| q)$

$$\underset{Y}{\mathbb{E}}[D(p_Y \| q)] = \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m)$$
$$+ \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m)$$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0,1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$,      $p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$,      $p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$\underset{Y}{\mathrm{E}}[D(p_Y \| q)] = \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m)$$
$$+ \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m)$$
$$= \Pr[Y = 0] \cdot \sum_i p_{0,i} \log \frac{p_{0,i}}{q_i} + \Pr[Y = 1] \cdot \sum_i p_{1,i} \log \frac{p_{1,i}}{q_i}$$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y=0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\qquad p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y=1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$, $\qquad p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$
\begin{aligned}
\mathop{\mathrm{E}}_{Y}[D(p_Y \| q)] &= \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m) \\
&\quad + \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m) \\
&= \Pr[Y = 0] \cdot \sum_i p_{0,i} \log \frac{p_{0,i}}{q_i} + \Pr[Y = 1] \cdot \sum_i p_{1,i} \log \frac{p_{1,i}}{q_i} \\
&= \Pr[Y = 0] \cdot \sum p_{0,i} \log p_{0,i} + \Pr[Y = 1] \cdot \sum p_{1,i} \log p_{1,i} \\
&\quad - \Pr[Y = 0] \cdot \sum p_{0,i} \log q_i - \Pr[Y = 1] \cdot \sum p_{1,i} \log q_i
\end{aligned}
$$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$,　　$p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$,　　$p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$
\begin{aligned}
\mathop{\mathrm{E}}_{Y}[D(p_Y \| q)] &= \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m) \\
&\quad + \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m) \\
&= \Pr[Y = 0] \cdot \sum_i p_{0,i} \log \frac{p_{0,i}}{q_i} + \Pr[Y = 1] \cdot \sum_i p_{1,i} \log \frac{p_{1,i}}{q_i} \\
&= \Pr[Y = 0] \cdot \sum p_{0,i} \log p_{0,i} + \Pr[Y = 1] \cdot \sum p_{1,i} \log p_{1,i} \\
&\quad - \Pr[Y = 0] \cdot \sum p_{0,i} \log q_i - \Pr[Y = 1] \cdot \sum p_{1,i} \log q_i \\
&= -H(X|Y) - \sum (\Pr[Y = 0] \cdot p_{0,i} + \Pr[Y = 1] \cdot p_{1,i}) \log q_i
\end{aligned}
$$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\quad p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$, $\quad p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$
\begin{aligned}
\mathbb{E}_Y[D(p_Y \| q)] &= \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m) \\
&\quad + \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m) \\
&= \Pr[Y = 0] \cdot \sum_i p_{0,i} \log \frac{p_{0,i}}{q_i} + \Pr[Y = 1] \cdot \sum_i p_{1,i} \log \frac{p_{1,i}}{q_i} \\
&= \Pr[Y = 0] \cdot \sum p_{0,i} \log p_{0,i} + \Pr[Y = 1] \cdot \sum p_{1,i} \log p_{1,i} \\
&\quad - \Pr[Y = 0] \cdot \sum p_{0,i} \log q_i - \Pr[Y = 1] \cdot \sum p_{1,i} \log q_i \\
&= -H(X|Y) - \sum (\Pr[Y = 0] \cdot p_{0,i} + \Pr[Y = 1] \cdot p_{1,i}) \log q_i \\
&= -H(X|Y) + H(X)
\end{aligned}
$$

## Mutual information as expected relative entropy

- Let $X \sim (q_1, \ldots, q_m)$ over $[m]$, and $Y$ be rv over $\{0, 1\}$
- $(X|Y = 0) \sim p_0 = (p_{0,1}, \ldots, p_{0,m})$, $\quad p_{0,i} = \Pr[X = i | Y = 0]$
- $(X|Y = 1) \sim p_1 = (p_{1,1}, \ldots, p_{1,m})$, $\quad p_{1,i} = \Pr[X = i | Y = 1]$
- If we learned $Y = j$, we gained $D(p_j \| q)$

$$
\begin{aligned}
\mathop{\mathrm{E}}_{Y}[D(p_Y \| q)] &= \Pr[Y = 0] \cdot D(p_{0,1}, \ldots, p_{0,m} \| q_1, \ldots, q_m) \\
&\quad + \Pr[Y = 1] \cdot D(p_{1,1}, \ldots, p_{1,m} \| q_1, \ldots, q_m) \\
&= \Pr[Y = 0] \cdot \sum_i p_{0,i} \log \frac{p_{0,i}}{q_i} + \Pr[Y = 1] \cdot \sum_i p_{1,i} \log \frac{p_{1,i}}{q_i} \\
&= \Pr[Y = 0] \cdot \sum p_{0,i} \log p_{0,i} + \Pr[Y = 1] \cdot \sum p_{1,i} \log p_{1,i} \\
&\quad - \Pr[Y = 0] \cdot \sum p_{0,i} \log q_i - \Pr[Y = 1] \cdot \sum p_{1,i} \log q_i \\
&= -H(X|Y) - \sum (\Pr[Y = 0] \cdot p_{0,i} + \Pr[Y = 1] \cdot p_{1,i}) \log q_i \\
&= -H(X|Y) + H(X) = I(X; Y)
\end{aligned}
$$

**Equivalent definition for mutual information**

## Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$

## Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$
- Interpretation

# Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p\|p_X p_Y)$
- Interpretation
- Proof:
  $$D(p\|p_X p_Y) = \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p_X(x) p_Y(y)}$$

# Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$
- Interpretation
- Proof:
  $$D(p \| p_X p_Y) = \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p_X(x) p_Y(y)}$$

**Equivalent definition for mutual information**

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$

- Interpretation

- Proof:
$$D(p \| p_X p_Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_X(x) p_Y(y)}$$
$$= \sum_{x,y} p(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)}$$

# Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$
- Interpretation
- Proof:

$$
\begin{aligned}
D(p \| p_X p_Y) &= \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p_X(x) p_Y(y)} \\
&= \sum_{x,y} p(x, y) \log \frac{p_{X|Y}(x|y)}{p_X(x)} \\
&= -\sum_{x,y} p(x, y) \log p_X(x) + \sum_{x,y} p(x, y) \log p_{X|Y}(x|y)
\end{aligned}
$$

## Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$
- Interpretation
- Proof:

$$D(p \| p_X p_Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_X(x) p_Y(y)}$$

$$= \sum_{x,y} p(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)}$$

$$= -\sum_{x,y} p(x,y) \log p_X(x) + \sum_{x,y} p(x,y) \log p_{X|Y}(x|y)$$

$$= H(X) + \sum_{y} p_Y(y) \sum_{x} p_{X|Y}(x|y) \log p_{X|Y}(x|y)$$

# Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$

- Interpretation

- Proof:
$$\begin{aligned}
D(p \| p_X p_Y) &= \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p_X(x) p_Y(y)} \\
&= \sum_{x,y} p(x, y) \log \frac{p_{X|Y}(x|y)}{p_X(x)} \\
&= -\sum_{x,y} p(x, y) \log p_X(x) + \sum_{x,y} p(x, y) \log p_{X|Y}(x|y) \\
&= H(X) + \sum_y p_Y(y) \sum_x p_{X|Y}(x|y) \log p_{X|Y}(x|y) \\
&= H(X) - H(X|Y)
\end{aligned}$$

## Equivalent definition for mutual information

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$
- Interpretation
- Proof:
$$D(p \| p_X p_Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_X(x) p_Y(y)}$$
$$= \sum_{x,y} p(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)}$$
$$= - \sum_{x,y} p(x,y) \log p_X(x) + \sum_{x,y} p(x,y) \log p_{X|Y}(x|y)$$
$$= H(X) + \sum_y p_Y(y) \sum_x p_{X|Y}(x|y) \log p_{X|Y}(x|y)$$
$$= H(X) - H(X|Y) = I(X; Y)$$

**Equivalent definition for mutual information**

- $(X, Y) \sim p$, then $I(X; Y) = D(p \| p_X p_Y)$

- Interpretation

- Proof:
$$D(p \| p_X p_Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p_X(x) p_Y(y)}$$
$$= \sum_{x,y} p(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)}$$
$$= -\sum_{x,y} p(x,y) \log p_X(x) + \sum_{x,y} p(x,y) \log p_{X|Y}(x|y)$$
$$= H(X) + \sum_y p_Y(y) \sum_x p_{X|Y}(x|y) \log p_{X|Y}(x|y)$$
$$= H(X) - H(X|Y) = I(X; Y)$$

- We will later see the relation between the above two facts.

Section 4

**Relation to Data Compression**

**Wrong code**

## Wrong code

**Theorem 2**

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$H(p) + D(p\|q) \leq \mathrm{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

## Wrong code

### Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$H(p) + D(p\|q) \leq \mathsf{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

▶ Recall that $H(q) \leq \mathsf{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.

## Wrong code

### Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. Then
$H(p) + D(p\|q) \leq \mathbb{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$*

- Recall that $H(q) \leq \mathbb{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

# Wrong code

## Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$$H(p) + D(p\|q) \leq \mathrm{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$$

- ▶ Recall that $H(q) \leq \mathrm{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- ▶ Proof of upperbound (upperbound is proved similarly)

# Wrong code

## Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$H(p) + D(p \| q) \leq \mathsf{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p \| q) + 1$

- Recall that $H(q) \leq \mathsf{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\mathsf{E}_{i \leftarrow p}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil$$

# Wrong code

**Theorem 2**

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with $\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. Then*
$$H(p) + D(p\|q) \leq \mathrm{E}_{i \leftarrow p}\left[\ell(i)\right] \leq H(p) + D(p\|q) + 1$$

- ▶ Recall that $H(q) \leq \mathrm{E}_{i \leftarrow q}\left[\ell(i)\right] \leq H(q) + 1$.
- ▶ Proof of upperbound (upperbound is proved similarly)

$$\underset{i \leftarrow p}{\mathrm{E}}\left[\ell(i)\right] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i (\log \frac{1}{q_i} + 1)$$

# Wrong code

## Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
*$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. Then*
*$H(p) + D(p\|q) \leq \mathsf{E}_{i\leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$*

- Recall that $H(q) \leq \mathsf{E}_{i\leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\mathsf{E}_{i\leftarrow p}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i(\log \frac{1}{q_i} + 1)$$

$$= 1 + \sum_i p_i(\log \frac{p_i}{q_i} \frac{1}{p_i})$$

# Wrong code

> **Theorem 2**
>
> Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with
> $\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. Then
> $H(p) + D(p\|q) \leq \mathsf{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

- Recall that $H(q) \leq \mathsf{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\mathsf{E}_{i \leftarrow p}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i (\log \frac{1}{q_i} + 1)$$

$$= 1 + \sum_i p_i (\log \frac{p_i}{q_i} \frac{1}{p_i}) = 1 + \sum_i p_i (\log \frac{p_i}{q_i}) + \sum_i p_i (\log \frac{1}{p_i})$$

# Wrong code

## Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$H(p) + D(p\|q) \leq E_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

- Recall that $H(q) \leq E_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\underset{i \leftarrow p}{E}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i(\log \frac{1}{q_i} + 1)$$

$$= 1 + \sum_i p_i(\log \frac{p_i}{q_i} \frac{1}{p_i}) = 1 + \sum_i p_i(\log \frac{p_i}{q_i}) + \sum_i p_i(\log \frac{1}{p_i})$$

$$= 1 + D(p\|q) + H(p)$$

# Wrong code

## Theorem 2

*Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
$\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
$H(p) + D(p\|q) \leq \mathrm{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

- ▶ Recall that $H(q) \leq \mathrm{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- ▶ Proof of upperbound (upperbound is proved similarly)

$$
\mathop{\mathrm{E}}_{i \leftarrow p}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i (\log \frac{1}{q_i} + 1)
$$
$$
= 1 + \sum_i p_i (\log \frac{p_i}{q_i} \frac{1}{p_i}) = 1 + \sum_i p_i (\log \frac{p_i}{q_i}) + \sum_i p_i (\log \frac{1}{p_i})
$$
$$
= 1 + D(p\|q) + H(p)
$$

- ▶ Can there be a (close) to optimal code for $q$ that is better for $p$?

## Wrong code

> **Theorem 2**
>
> *Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
> $\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
> $H(p) + D(p\|q) \leq \mathsf{E}_{i \leftarrow p}[\ell(i)] \leq H(p) + D(p\|q) + 1$

- Recall that $H(q) \leq \mathsf{E}_{i \leftarrow q}[\ell(i)] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\mathsf{E}_{i \leftarrow p}[\ell(i)] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i (\log \frac{1}{q_i} + 1)$$

$$= 1 + \sum_i p_i (\log \frac{p_i}{q_i} \frac{1}{p_i}) = 1 + \sum_i p_i (\log \frac{p_i}{q_i}) + \sum_i p_i (\log \frac{1}{p_i})$$

$$= 1 + D(p\|q) + H(p)$$

- Can there be a (close) to optimal code for $q$ that is better for $p$?

# Wrong code

> **Theorem 2**
>
> *Let $p$ and $q$ be distributions over $[m]$, and let $C$ be code with*
> $\ell(i) = |C(i)| = \left\lceil \log \frac{1}{q_i} \right\rceil$. *Then*
> $H(p) + D(p\|q) \leq \mathsf{E}_{i \leftarrow p}\left[\ell(i)\right] \leq H(p) + D(p\|q) + 1$

- Recall that $H(q) \leq \mathsf{E}_{i \leftarrow q}\left[\ell(i)\right] \leq H(q) + 1$.
- Proof of upperbound (upperbound is proved similarly)

$$\mathsf{E}_{i \leftarrow p}\left[\ell(i)\right] = \sum_i p_i \left\lceil \log \frac{1}{q_i} \right\rceil < \sum_i p_i (\log \frac{1}{q_i} + 1)$$

$$= 1 + \sum_i p_i (\log \frac{p_i}{q_i} \frac{1}{p_i}) = 1 + \sum_i p_i (\log \frac{p_i}{q_i}) + \sum_i p_i (\log \frac{1}{p_i})$$

$$= 1 + D(p\|q) + H(p)$$

- Can there be a (close) to optimal code for $q$ that is better for $p$? HW

Section 5

**Conditional Relative Entropy**

**Conditional relative entropy**

## Conditional relative entropy

**Definition 3**

For two distributions *p* and *q* over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

## Conditional relative entropy

**Definition 3**

For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \mathop{\mathrm{E}}_{(X,Y) \sim p(x,y)} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

## Conditional relative entropy

**Definition 3**

For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \underset{(X,Y) \sim p(x,y)}{\mathsf{E}} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

▶ Let $(X_p, Y_p) \sim p$ and $(X_q, Y_q) \sim q$, then
  $D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \mathsf{E}_{x \leftarrow X_p} \left[ D(Y_q|_{x_p=x} \| Y_q|_{x_q=x}) \right]$

## Conditional relative entropy

**Definition 3**

For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \underset{(X,Y) \sim p(x,y)}{\mathsf{E}} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

▶ Let $(X_p, Y_p) \sim p$ and $(X_q, Y_q) \sim q$, then
$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \mathsf{E}_{x \leftarrow X_p} \left[ D(Y_q|_{x_p=x} \| Y_q|_{x_q=x}) \right]$

▶ Example:

$$p = \begin{array}{c|c|c} {}_{X}\diagdown^{Y} & 0 & 1 \\ \hline 0 & \frac{1}{8} & \frac{1}{8} \\ \hline 1 & \frac{1}{4} & \frac{1}{2} \end{array}$$

$$q = \begin{array}{c|c|c} {}_{X}\diagdown^{Y} & 0 & 1 \\ \hline 0 & \frac{1}{8} & \frac{1}{4} \\ \hline 1 & \frac{1}{2} & \frac{1}{8} \end{array}$$

# Conditional relative entropy

**Definition 3**

For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \mathop{\mathrm{E}}_{(X,Y) \sim p(x,y)} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

▶ Let $(X_p, Y_p) \sim p$ and $(X_q, Y_q) \sim q$, then
$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \mathrm{E}_{x \leftarrow X_p} \left[ D(Y_q|_{x_p=x} \| Y_q|_{x_q=x}) \right]$

▶ Example:

$$p = \begin{array}{c|c|c} {}_x\backslash^Y & 0 & 1 \\ \hline 0 & \frac{1}{8} & \frac{1}{8} \\ \hline 1 & \frac{1}{4} & \frac{1}{2} \end{array} \qquad q = \begin{array}{c|c|c} {}_x\backslash^Y & 0 & 1 \\ \hline 0 & \frac{1}{8} & \frac{1}{4} \\ \hline 1 & \frac{1}{2} & \frac{1}{8} \end{array}$$

# Conditional relative entropy

> **Definition 3**
>
> For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:
>
> $$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$
>
> $$= \mathop{\mathrm{E}}_{(X,Y) \sim p(x,y)} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

- Let $(X_p, Y_p) \sim p$ and $(X_q, Y_q) \sim q$, then
  $D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \mathrm{E}_{x \leftarrow X_p} \left[ D(Y_q|_{x_p=x} \| Y_q|_{x_q=x}) \right]$

- Example:

  $p = $

  | $x$\\$^Y$ | 0 | 1 |
  |-----------|---|---|
  | 0 | $\frac{1}{8}$ | $\frac{1}{8}$ |
  | 1 | $\frac{1}{4}$ | $\frac{1}{2}$ |

  $q = $

  | $x$\\$^Y$ | 0 | 1 |
  |-----------|---|---|
  | 0 | $\frac{1}{8}$ | $\frac{1}{4}$ |
  | 1 | $\frac{1}{2}$ | $\frac{1}{8}$ |

  $$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \frac{1}{4} \cdot D((\frac{1}{2}, \frac{1}{2}) \| (\frac{1}{3}, \frac{2}{3})) + \frac{3}{4} \cdot D((\frac{1}{3}, \frac{2}{3}) \| (\frac{4}{5}, \frac{1}{5}))$$

## Conditional relative entropy

**Definition 3**

For two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$:

$$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) := \sum_{x \in \mathcal{X}} p_{\mathcal{X}}(x) \cdot \sum_{y \in \mathcal{Y}} p_{\mathcal{Y}|\mathcal{X}}(y|x) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \operatorname*{E}_{(X,Y) \sim p(x,y)} \left[ \log \frac{p_{\mathcal{Y}|\mathcal{X}}(Y|X)}{q_{\mathcal{Y}|\mathcal{X}}(Y|X)} \right]$$

- Let $(X_p, Y_p) \sim p$ and $(X_q, Y_q) \sim q$, then
  $D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = E_{x \leftarrow X_p} \left[ D(Y_q|_{x_p=x} \| Y_q|_{x_q=x}) \right]$

- Example:

  $p = $
  | $x$\\$Y$ | 0 | 1 |
  |---|---|---|
  | 0 | $\frac{1}{8}$ | $\frac{1}{8}$ |
  | 1 | $\frac{1}{4}$ | $\frac{1}{2}$ |

  $q = $
  | $x$\\$Y$ | 0 | 1 |
  |---|---|---|
  | 0 | $\frac{1}{8}$ | $\frac{1}{4}$ |
  | 1 | $\frac{1}{2}$ | $\frac{1}{8}$ |

  $$D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}}) = \frac{1}{4} \cdot D((\frac{1}{2}, \frac{1}{2}) \| (\frac{1}{3}, \frac{2}{3})) + \frac{3}{4} \cdot D((\frac{1}{3}, \frac{2}{3}) \| (\frac{4}{5}, \frac{1}{5}))$$

  $$= \dots$$

# Chain rule

## Chain rule

**Claim 4**

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

# Chain rule

## Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_\mathcal{X}\|q_\mathcal{X}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

- Proof:

$$D(p\|q) = \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

# Chain rule

## Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

- ▶ Proof:

$$D(p\|q) = \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

# Chain rule

**Claim 4**

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

▶ Proof:

$$D(p\|q) = \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

$$= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x)q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

# Chain rule

### Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

- Proof:

$$D(p\|q) = \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

$$= \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x)q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

$$= \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)}$$

# Chain rule

## Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p \| q) = D(p_{\mathcal{X}} \| q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}})$

- Proof:

$$
\begin{aligned}
D(p \| q) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x) p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x) q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= D(p_{\mathcal{X}} \| q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}} \| q_{\mathcal{Y}|\mathcal{X}})
\end{aligned}
$$

# Chain rule

## Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

- Proof:

$$
\begin{aligned}
D(p\|q) &= \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)} \\
&= \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x)q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y)\in\mathcal{X}\times\mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})
\end{aligned}
$$

- It follows that for $(X;Y) \sim p$: $I(X,Y) = D(p\|p_X p_Y)$

## Chain rule

**Claim 4**

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

▶ Proof:

$$
\begin{aligned}
D(p\|q) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x) p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x) q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})
\end{aligned}
$$

▶ It follows that for $(X;Y) \sim p$: $I(X,Y) = D(p\|p_X p_Y)$

# Chain rule

<div>

**Claim 4**

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

</div>

▶ Proof:

$$
\begin{aligned}
D(p\|q) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x)q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})
\end{aligned}
$$

▶ It follows that for $(X;Y) \sim p$: $I(X,Y) = D(p\|p_X p_Y) = D(p_X\|p_X) + \mathrm{E}_{x \leftarrow X}\left[D(p_{Y|_{X=x}}, p_Y)\right]$

## Chain rule

### Claim 4

For any two distributions $p$ and $q$ over $\mathcal{X} \times \mathcal{Y}$, it holds that
$D(p\|q) = D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})$

▶ Proof:

$$
\begin{aligned}
D(p\|q) &= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{q(x,y)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x) p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{X}}(x) q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{X}}(x)}{q_{\mathcal{X}}(x)} + \sum_{(x,y) \in \mathcal{X} \times \mathcal{Y}} p(x,y) \log \frac{p_{\mathcal{Y}|\mathcal{X}}(y|x)}{q_{\mathcal{Y}|\mathcal{X}}(y|x)} \\
&= D(p_{\mathcal{X}}\|q_{\mathcal{X}}) + D(p_{\mathcal{Y}|\mathcal{X}}\|q_{\mathcal{Y}|\mathcal{X}})
\end{aligned}
$$

▶ It follows that for $(X;Y) \sim p$: $I(X,Y) = D(p\|p_X p_Y) = $
$D(p_X\|p_X) + \mathrm{E}_{x \leftarrow X}\left[D(p_{Y|X=x}, p_Y)\right] = \mathrm{E}_{x \leftarrow X}\left[D(p_{Y|X=x}, p_Y)\right]$

Section 6

**Data-processing inequality**

# Data-processing inequality

# Data-processing inequality

**Claim 5**

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

## Data-processing inequality

**Claim 5**

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

▸ Analogues to $H(X) \geq H(f(X))$

# Data-processing inequality

**Claim 5**

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

- Analogues to $H(X) \geq H(f(X))$
- Proof:

## Data-processing inequality

### Claim 5

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

- Analogues to $H(X) \geq H(f(X))$
- Proof:
- $D(X, f(X)\|Y, f(Y)) = D(X\|Y)$

## Data-processing inequality

### Claim 5

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

- Analogues to $H(X) \geq H(f(X))$
- Proof:
- $D(X, f(X)\|Y, f(Y)) = D(X\|Y)$
- $D(X, f(X)\|Y, f(Y)) = D(f(X)\|f(Y)) + \mathsf{E}_{z \leftarrow f(X)}\left[D(X|_{f(X)=z}\|Y|_{f(X)=z}))\right] \geq D(f(X)\|f(Y))$

# Data-processing inequality

## Claim 5

For any rv's $X$ and $Y$ and function $f$, it holds that $D(f(X)\|f(Y)) \leq D(X\|Y)$.

- Analogues to $H(X) \geq H(f(X))$
- Proof:
- $D(X, f(X)\|Y, f(Y)) = D(X\|Y)$
- $D(X, f(X)\|Y, f(Y)) = D(f(X)\|f(Y)) + \mathsf{E}_{z \leftarrow f(X)}\left[D(X|_{f(X)=z}\|Y|_{f(X)=z}))\right] \geq D(f(X)\|f(Y))$

- Hence, $D(f(X)\|f(Y)) \leq D(X\|Y)$.

# Section 7

## **Relation to Statistical Distance**

**Relation to statistical distance**

## Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

## Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

## Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$
- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality
- However,

## Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

- However,

**Theorem 6**

$\mathrm{SD}(p, q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

# Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

- However,

**Theorem 6**

$\mathrm{SD}(p, q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

# Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

- However,

**Theorem 6**

$\mathrm{SD}(p, q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

- Corollary: For rv $X$ over $[m]$ with $H(X) \geq m - \varepsilon$, it holds that
  $\mathrm{SD}(X, \sim [m]) \leq \sqrt{\frac{\ln 2}{2} \cdot (m - H(X))} = \sqrt{\frac{\ln 2}{2} \cdot \varepsilon}$

# Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$
- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality
- However,

> **Theorem 6**
>
> $\mathrm{SD}(p,q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

- Corollary: For rv $X$ over $[m]$ with $H(X) \geq m - \varepsilon$, it holds that
  $$\mathrm{SD}(X, \sim [m]) \leq \sqrt{\frac{\ln 2}{2} \cdot (m - H(X))} = \sqrt{\frac{\ln 2}{2} \cdot \varepsilon}$$
- Other direction is incorrect: $\mathrm{SD}(p,q)$ might be small but $D(p\|q) = \infty$

# Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

- However,

> **Theorem 6**
>
> $\mathsf{SD}(p, q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

- Corollary: For rv $X$ over $[m]$ with $H(X) \geq m - \varepsilon$, it holds that

  $\mathsf{SD}(X, \sim [m]) \leq \sqrt{\frac{\ln 2}{2} \cdot (m - H(X))} = \sqrt{\frac{\ln 2}{2} \cdot \varepsilon}$

- Other direction is incorrect: $\mathsf{SD}(p, q)$ might be small but $D(p\|q) = \infty$

- Does $\mathsf{SD}(p, \sim [m])$ being small imply $D(p\| \sim [m]) = \log m - H(p)$ is small?

# Relation to statistical distance

- $D(p\|q)$ is used many time to measure the distance from $p$ to $q$

- It is not a distance in the mathematical sense: $D(p\|q) \neq D(q\|p)$ and no triangle inequality

- However,

> **Theorem 6**
>
> $\mathrm{SD}(p, q) \leq \sqrt{\frac{\ln 2}{2} \cdot D(p\|q)}$

- Corollary: For rv $X$ over $[m]$ with $H(X) \geq m - \varepsilon$, it holds that
  $$\mathrm{SD}(X, \sim [m]) \leq \sqrt{\frac{\ln 2}{2} \cdot (m - H(X))} = \sqrt{\frac{\ln 2}{2} \cdot \varepsilon}$$

- Other direction is incorrect: $\mathrm{SD}(p, q)$ might be small but $D(p\|q) = \infty$

- Does $\mathrm{SD}(p, \sim [m])$ being small imply $D(p\| \sim [m]) = \log m - H(p)$ is small?

  HW

# Proving Thm 6, boolean case

# Proving Thm 6, boolean case

- Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

# Proving Thm 6, boolean case

▶ Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

▶ $SD(p, q) = \alpha - \beta$

# Proving Thm 6, boolean case

▶ Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

▶ $\mathrm{SD}(p, q) = \alpha - \beta$

▶ We will show that
$D(p\|q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2\ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2}\mathrm{SD}(p, q)^2$

# Proving Thm 6, boolean case

▶ Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

▶ $SD(p, q) = \alpha - \beta$

▶ We will show that
$D(p \| q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2\ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2}SD(p, q)^2$

▶ Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} - \frac{4}{2\ln 2}(\alpha - \beta)^2$

# Proving Thm 6, boolean case

- ▶ Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$
- ▶ $SD(p, q) = \alpha - \beta$
- ▶ We will show that
  $D(p\|q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2 \ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2} SD(p, q)^2$
- ▶ Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} - \frac{4}{2 \ln 2}(\alpha - \beta)^2$

$$\frac{\partial g(x, y)}{\partial y} = -\frac{x}{y \ln 2} + \frac{1 - x}{(1 - y) \ln 2} - \frac{4}{2 \ln 2} 2(y - x)$$

# Proving Thm 6, boolean case

- Let $p = (\alpha, 1-\alpha)$ and $q = (\beta, 1-\beta)$ and assume $\alpha \geq \beta$

- $\mathrm{SD}(p, q) = \alpha - \beta$

- We will show that
  $D(p\|q) = \alpha \log \frac{\alpha}{\beta} + (1-\alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2\ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2}\mathrm{SD}(p, q)^2$

- Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1-\alpha) \log \frac{1-\alpha}{1-\beta} - \frac{4}{2\ln 2}(\alpha - \beta)^2$

$$\frac{\partial g(x, y)}{\partial y} = -\frac{x}{y \ln 2} + \frac{1-x}{(1-y)\ln 2} - \frac{4}{2\ln 2}2(y - x)$$

# Proving Thm 6, boolean case

- Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

- $\mathrm{SD}(p, q) = \alpha - \beta$

- We will show that
  $D(p \| q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1 - \alpha}{1 - \beta} \geq \frac{4}{2 \ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2}\mathrm{SD}(p, q)^2$

- Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1 - \alpha}{1 - \beta} - \frac{4}{2 \ln 2}(\alpha - \beta)^2$

$$\frac{\partial g(x, y)}{\partial y} = -\frac{x}{y \ln 2} + \frac{1 - x}{(1 - y) \ln 2} - \frac{4}{2 \ln 2} 2(y - x)$$

$$= \frac{y - x}{y(1 - y) \ln 2} - \frac{4}{\ln 2}(y - x)$$

## Proving Thm 6, boolean case

- Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

- $\mathrm{SD}(p, q) = \alpha - \beta$

- We will show that
  $D(p\|q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2 \ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2}\mathrm{SD}(p, q)^2$

- Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} - \frac{4}{2 \ln 2}(\alpha - \beta)^2$

$$\frac{\partial g(x, y)}{\partial y} = -\frac{x}{y \ln 2} + \frac{1 - x}{(1 - y) \ln 2} - \frac{4}{2 \ln 2} 2(y - x)$$

$$= \frac{y - x}{y(1 - y) \ln 2} - \frac{4}{\ln 2}(y - x)$$

- Since $y(1 - y) \leq \frac{1}{4}$, $\frac{\partial g(x, y)}{\partial y} \leq 0$ for $y < x$.

## Proving Thm 6, boolean case

▶ Let $p = (\alpha, 1 - \alpha)$ and $q = (\beta, 1 - \beta)$ and assume $\alpha \geq \beta$

▶ $SD(p, q) = \alpha - \beta$

▶ We will show that
  $D(p\|q) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} \geq \frac{4}{2 \ln 2}(\alpha - \beta)^2 = \frac{2}{\ln 2} SD(p, q)^2$

▶ Let $g(\alpha, \beta) = \alpha \log \frac{\alpha}{\beta} + (1 - \alpha) \log \frac{1-\alpha}{1-\beta} - \frac{4}{2 \ln 2}(\alpha - \beta)^2$

$$\frac{\partial g(x, y)}{\partial y} = -\frac{x}{y \ln 2} + \frac{1 - x}{(1 - y) \ln 2} - \frac{4}{2 \ln 2} 2(y - x)$$

$$= \frac{y - x}{y(1 - y) \ln 2} - \frac{4}{\ln 2}(y - x)$$

▶ Since $y(1 - y) \leq \frac{1}{4}$, $\frac{\partial g(x, y)}{\partial y} \leq 0$ for $y < x$.

▶ Since $g(x, x) = 0$, $g(x, y) \geq 0$ for $y < x$. $\square$

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$

# Proving Thm 6, general case

- Let $\mathcal{U} = \mathrm{Supp}(p) \cup \mathrm{Supp}(q)$
- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

# Proving Thm 6, general case

- Let $\mathcal{U} = \mathrm{Supp}(p) \cup \mathrm{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

- $\mathrm{SD}(p, q) = \mathrm{Pr}_p\left[\mathcal{S}\right] - \mathrm{Pr}_q\left[\mathcal{S}\right]$    (by homework)

# Proving Thm 6, general case

- Let $\mathcal{U} = \mathrm{Supp}(p) \cup \mathrm{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

- $\mathrm{SD}(p, q) = \mathrm{Pr}_p[\mathcal{S}] - \mathrm{Pr}_q[\mathcal{S}]$    (by homework)

- Let $P \sim p$, and let the indicator $\hat{P}$ be $1$ iff $P \in \mathcal{S}$.

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$ \qquad (by homework)

- Let $P \sim p$, and let the indicator $\hat{P}$ be $1$ iff $P \in \mathcal{S}$.

- Let $Q \sim q$, and let the indicator $\hat{Q}$ be $1$ iff $Q \in \mathcal{S}$.

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$     (by homework)

- Let $P \sim p$, and let the indicator $\hat{P}$ be $1$ iff $P \in \mathcal{S}$.

- Let $Q \sim q$, and let the indicator $\hat{Q}$ be $1$ iff $Q \in \mathcal{S}$.

- $\text{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \text{SD}(p, q)$

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} : p(u) > q(u)\}$

- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$     (by homework)

- Let $P \sim p$, and let the indicator $\hat{P}$ be 1 iff $P \in \mathcal{S}$.

- Let $Q \sim q$, and let the indicator $\hat{Q}$ be 1 iff $Q \in \mathcal{S}$.

- $\text{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \text{SD}(p, q)$

  $D(p \| q) \geq D(\hat{P} \| \hat{Q})$

## Proving Thm 6, general case

- Let $\mathcal{U} = \mathrm{Supp}(p) \cup \mathrm{Supp}(q)$
- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$
- $\mathrm{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$     (by homework)
- Let $P \sim p$, and let the indicator $\hat{P}$ be 1 iff $P \in \mathcal{S}$.
- Let $Q \sim q$, and let the indicator $\hat{Q}$ be 1 iff $Q \in \mathcal{S}$.
- $\mathrm{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \mathrm{SD}(p, q)$

$$D(p \| q) \geq D(\hat{P} \| \hat{Q})$$

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$

- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$

- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$     (by homework)

- Let $P \sim p$, and let the indicator $\hat{P}$ be 1 iff $P \in \mathcal{S}$.

- Let $Q \sim q$, and let the indicator $\hat{Q}$ be 1 iff $Q \in \mathcal{S}$.

- $\text{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \text{SD}(p, q)$

$$D(p\|q) \geq D(\hat{P}\|\hat{Q})$$      (data-proccessing inequality)

# Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$
- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$
- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$   (by homework)
- Let $P \sim p$, and let the indicator $\hat{P}$ be 1 iff $P \in \mathcal{S}$.
- Let $Q \sim q$, and let the indicator $\hat{Q}$ be 1 iff $Q \in \mathcal{S}$.
- $\text{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \text{SD}(p, q)$

$$
\begin{aligned}
D(p \| q) &\geq D(\hat{P} \| \hat{Q}) \qquad \text{(data-proccessing inequality)} \\
&\geq \frac{2}{\ln 2} \cdot \text{SD}(\hat{P}, \hat{Q})^2
\end{aligned}
$$

## Proving Thm 6, general case

- ▶ Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$
- ▶ Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$
- ▶ $\text{SD}(p, q) = \text{Pr}_p[\mathcal{S}] - \text{Pr}_q[\mathcal{S}]$    (by homework)
- ▶ Let $P \sim p$, and let the indicator $\hat{P}$ be $1$ iff $P \in \mathcal{S}$.
- ▶ Let $Q \sim q$, and let the indicator $\hat{Q}$ be $1$ iff $Q \in \mathcal{S}$.
- ▶ $\text{SD}(\hat{P}, \hat{Q}) = \text{Pr}[P \in \mathcal{S}] - \text{Pr}[Q \in \mathcal{S}] = \text{SD}(p, q)$

$$
\begin{aligned}
D(p\|q) &\geq D(\hat{P}\|\hat{Q}) && \text{(data-proccessing inequality)} \\
&\geq \frac{2}{\ln 2} \cdot \text{SD}(\hat{P}, \hat{Q})^2 && \text{(the Boolean case)}
\end{aligned}
$$

## Proving Thm 6, general case

- Let $\mathcal{U} = \text{Supp}(p) \cup \text{Supp}(q)$
- Let $\mathcal{S} = \{u \in \mathcal{U} \colon p(u) > q(u)\}$
- $\text{SD}(p, q) = \Pr_p[\mathcal{S}] - \Pr_q[\mathcal{S}]$     (by homework)
- Let $P \sim p$, and let the indicator $\hat{P}$ be $1$ iff $P \in \mathcal{S}$.
- Let $Q \sim q$, and let the indicator $\hat{Q}$ be $1$ iff $Q \in \mathcal{S}$.
- $\text{SD}(\hat{P}, \hat{Q}) = \Pr[P \in \mathcal{S}] - \Pr[Q \in \mathcal{S}] = \text{SD}(p, q)$

$$
\begin{aligned}
D(p\|q) &\geq D(\hat{P}\|\hat{Q}) && \text{(data-proccessing inequality)} \\
&\geq \frac{2}{\ln 2} \cdot \text{SD}(\hat{P}, \hat{Q})^2 && \text{(the Boolean case)} \\
&= \frac{2}{\ln 2} \cdot \text{SD}(p, q)^2. \ \square
\end{aligned}
$$

Section 8

**Conditioned Distributions**

# Main theorem

## Main theorem

**Theorem 7**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.*

## Main theorem

**Theorem 7**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.*

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines.

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$D(Y \| X) = \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})}$$

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$D(Y \| X) = \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})} = \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(y_1)} \frac{Y_2(y_2)}{X_2(y_2)} \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)}$$

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$
\begin{aligned}
D(Y \| X) &= \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})} = \sum_{\mathbf{y}=(y_1,y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(y_1)} \frac{Y_2(y_2)}{X_2(y_2)} \frac{Y(\mathbf{y})}{Y_1(y_1)Y_2(y_2)} \\
&= \sum_{\mathbf{y}=(y_1,y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(x_1)} + \sum_{\mathbf{y}=(y_1,y_2)} Y(\mathbf{y}) \log \frac{Y_2(y_2)}{X_2(x_1)} \\
&\quad + \sum_{\mathbf{y}=(y_1,y_2)} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{Y_1(y_1)Y_2(y_2)}
\end{aligned}
$$

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$D(Y \| X) = \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})} = \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(y_1)} \frac{Y_2(y_2)}{X_2(y_2)} \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)}$$

$$= \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(x_1)} + \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_2(y_2)}{X_2(x_1)}$$

$$+ \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)}$$

$$= D(Y_1 \| X_1) + D(Y_2 \| X_2) + I(Y_1; Y_2)$$

## Main theorem

**Theorem 7**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$
\begin{aligned}
D(Y \| X) &= \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})} = \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(y_1)} \frac{Y_2(y_2)}{X_2(y_2)} \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)} \\
&= \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(x_1)} + \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_2(y_2)}{X_2(x_1)} \\
&\quad + \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)} \\
&= D(Y_1 \| X_1) + D(Y_2 \| X_2) + I(Y_1; Y_2) \geq D(Y_1 \| X_1) + D(Y_2 \| X_2)
\end{aligned}
$$

## Main theorem

**Theorem 7**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{U}$, and let $Y = (Y_1, \ldots, Y_k)$ be rv over $\mathcal{U}^k$. Then $\sum_{j=1}^{k} D(Y_j \| X_j) \leq D(Y \| (X_1, \ldots, X_k))$.*

For rv $Z$, let $Z(z) = \Pr[Z = z]$.

We prove for $k = 2$, general case follows similar lines. Let $X = (X_1, X_2)$

$$D(Y \| X) = \sum_{\mathbf{y} \in \mathcal{U}^2} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{X(\mathbf{y})} = \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(y_1)} \frac{Y_2(y_2)}{X_2(y_2)} \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)}$$

$$= \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_1(y_1)}{X_1(x_1)} + \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y_2(y_2)}{X_2(x_1)}$$

$$+ \sum_{\mathbf{y} = (y_1, y_2)} Y(\mathbf{y}) \log \frac{Y(\mathbf{y})}{Y_1(y_1) Y_2(y_2)}$$

$$= D(Y_1 \| X_1) + D(Y_2 \| X_2) + I(Y_1; Y_2) \geq D(Y_1 \| X_1) + D(Y_2 \| X_2)$$

**Conditioning distributions, relative entropy case**

## Conditioning distributions, relative entropy case

### Theorem 8

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W) \| X_j) \leq D((X|_W) \| X) \leq \log \frac{1}{\Pr[W]}$.*

## Conditioning distributions, relative entropy case

**Theorem 8**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X)$$

## Conditioning distributions, relative entropy case

**Theorem 8**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \qquad\qquad \text{(Thm 7)}$$

## Conditioning distributions, relative entropy case

**Theorem 8**

Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \leq \log \frac{1}{\Pr[W]}$.

$$\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \qquad \text{(Thm 7)}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_W)(\mathbf{x}) \log \frac{(X|_W)(\mathbf{x})}{X(\mathbf{x})}$$

## Conditioning distributions, relative entropy case

### Theorem 8

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W) \| X_j) \leq D((X|_W) \| X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_W) \| X_j) \leq D((X|_W) \| X) \qquad \text{(Thm 7)}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_W)(\mathbf{x}) \log \frac{(X|_W)(\mathbf{x})}{X(\mathbf{x})}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_W)(\mathbf{x}) \log \frac{\Pr[W|X = \mathbf{x}])}{\Pr[W]}$$

## Conditioning distributions, relative entropy case

**Theorem 8**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_W)\|X_j) \leq D((X|_W)\|X) \qquad \text{(Thm 7)}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_W)(\mathbf{x}) \log \frac{(X|_W)(\mathbf{x})}{X(\mathbf{x})}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_W)(\mathbf{x}) \log \frac{\Pr[W|X = \mathbf{x}])}{\Pr[W]} \qquad \text{(Bayes)}$$

## Conditioning distributions, relative entropy case

**Theorem 8**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_w)\|X_j) \leq D((X|_w)\|X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_w)\|X_j) \leq D((X|_w)\|X) \qquad \text{(Thm 7)}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \frac{(X|_w)(\mathbf{x})}{X(\mathbf{x})}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \frac{\Pr[W|X = \mathbf{x}])}{\Pr[W]} \qquad \text{(Bayes)}$$

$$= \log \frac{1}{\Pr[W]} + \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \Pr[W|X = \mathbf{x}])$$

## Conditioning distributions, relative entropy case

**Theorem 8**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$, let $X = (X_1, \ldots, X_k)$ and let $W$ be an event (i.e., Boolean rv). Then $\sum_{j=1}^{k} D((X_j|_w)\|X_j) \leq D((X|_w)\|X) \leq \log \frac{1}{\Pr[W]}$.*

$$\sum_{j=1}^{k} D((X_j|_w)\|X_j) \leq D((X|_w)\|X) \qquad \text{(Thm 7)}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \frac{(X|_w)(\mathbf{x})}{X(\mathbf{x})}$$

$$= \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \frac{\Pr[W|X = \mathbf{x}])}{\Pr[W]} \qquad \text{(Bayes)}$$

$$= \log \frac{1}{\Pr[W]} + \sum_{\mathbf{x} \in \mathcal{X}^k} (X|_w)(\mathbf{x}) \log \Pr[W|X = \mathbf{x}])$$

$$\leq \log \frac{1}{\Pr[W]}$$

**Conditioning distributions, statistical distance case**

# Conditioning distributions, statistical distance case

## Theorem 9

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$ and let $W$ be an event. Then $\sum_{j=1}^{k} \mathsf{SD}((X_j|_W), X_j)^2 \leq \log \frac{1}{\Pr[W]}$.*

# Conditioning distributions, statistical distance case

**Theorem 9**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$ and let $W$ be an event. Then $\sum_{j=1}^{k} \mathsf{SD}((X_j|_W), X_j)^2 \leq \log \frac{1}{\Pr[W]}$.*

Proof: follows by Thm 6, and Thm 7.□

# Conditioning distributions, statistical distance case

## Theorem 9

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$ and let $W$ be an event. Then*
$\sum_{j=1}^{k} \mathsf{SD}((X_j|_W), X_j)^2 \leq \log \frac{1}{\Pr[W]}$.

Proof: follows by Thm 6, and Thm 7.□

Using $(\sum_{j=1}^{k} a_i)^2 \leq k \cdot \sum_{i=1}^{k} a_i^2$, it follows that

## Corollary 10

$\sum_{j=1}^{k} \mathsf{SD}((X_j|_W), X_j) \leq \sqrt{k \log(\frac{1}{\Pr[W]})}$, *and*

$\mathsf{E}_{j \leftarrow k} \, \mathsf{SD}((X_j|_W), X_j) \leq \sqrt{\frac{1}{k} \log(\frac{1}{\Pr[W]})}$

## Conditioning distributions, statistical distance case

---

**Theorem 9**

*Let $X_1, \ldots, X_k$ be iid over $\mathcal{X}$ and let $W$ be an event. Then*
$\sum_{j=1}^{k} \mathsf{SD}((X_j|_w), X_j)^2 \leq \log \frac{1}{\Pr[W]}$.

---

Proof: follows by Thm 6, and Thm 7.$\square$

Using $(\sum_{j=1}^{k} a_i)^2 \leq k \cdot \sum_{i=1}^{k} a_i^2$, it follows that

---

**Corollary 10**

$\sum_{j=1}^{k} \mathsf{SD}((X_j|_w), X_j) \leq \sqrt{k \log(\frac{1}{\Pr[W]})}$, *and*

$\mathsf{E}_{j \leftarrow k} \mathsf{SD}((X_j|_w), X_j) \leq \sqrt{\frac{1}{k} \log(\frac{1}{\Pr[W]})}$

---

Extraction

# Numerical example

## Numerical example

- Let $X = (X_1, \ldots, X_k) \leftarrow \{0, 1\}^{40}$ and let $f \colon \{0, 1\}^{40} \mapsto 0$ be such that $\Pr[f(X) = 0] = 2^{-10}$.

## Numerical example

- Let $X = (X_1, \ldots, X_k) \leftarrow \{0, 1\}^{40}$ and let $f \colon \{0, 1\}^{40} \mapsto 0$ be such that $\Pr[f(X) = 0] = 2^{-10}$.

- $\mathsf{E}_{j \leftarrow [40]} \mathsf{SD}((X_j|_{f(X)=0}), \sim \{0, 1\}) \leq \sqrt{\frac{1}{40} \cdot 10} = \frac{1}{2}$

## Numerical example

- Let $X = (X_1, \ldots, X_k) \leftarrow \{0, 1\}^{40}$ and let $f \colon \{0, 1\}^{40} \mapsto 0$ be such that $\Pr[f(X) = 0] = 2^{-10}$.

- $\mathsf{E}_{j \leftarrow [40]} \mathsf{SD}((X_j|_{f(X)=0}), \sim \{0, 1\}) \leq \sqrt{\frac{1}{40} \cdot 10} = \frac{1}{2}$

- Typical bits are not too biassed, even when conditioning on a very unlikely event.

# Extension

## Extension

**Theorem 11**

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively. Let $W$ be an event and assume that the $X_i$'s are iid conditioned on $T$. Then $\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T)) \leq \log \frac{1}{\Pr[W]} + \log |\text{Supp}(V|_W)|$, where $X_j'(t)$ is distributed according to $X_j|_{T=t}$.

## Extension

### Theorem 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively. Let $W$ be an event and assume that the $X_i$'s are iid conditioned on $T$. Then
$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T)) \leq \log \frac{1}{\Pr[W]} + \log |\mathsf{Supp}(V|_W)|$,
where $X_j'(t)$ is distributed according to $X_j|_{T=t}$.

Interpretation.

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X'_j(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X'_j(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X'_j(T))$$
$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \left[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t}) \right]$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X'_j(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X'_j(T))$$
$$= \underset{(t,v) \leftarrow (TV)|_W}{\mathbb{E}} \Big[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t})) \Big] \qquad \text{(chain rule)}$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))$$

$$= \underset{(t,v) \leftarrow (TV)|_W}{\mathsf{E}} \left[ \sum_{j=1}^{k} D(X_j|_{W, V=v, T=t} \| (X_j|_{T=t}) \right] \qquad \text{(chain rule)}$$

$$= \underset{(t,v) \leftarrow (TV)|_W}{\mathsf{E}} \left[ \sum D((X_j|\underbrace{W, V = v}_{w'})|_{T=t} \| (X_j|T = t)) \right]$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))$$

$$= \underset{(t,v) \leftarrow (TV)|_W}{\mathrm{E}} \Big[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t}) \Big] \qquad \text{(chain rule)}$$

$$= \underset{(t,v) \leftarrow (TV)|_W}{\mathrm{E}} \Big[ \sum D((X_j|\underbrace{W,\, V=v}_{w'})|_{T=t} \| (X_j|T=t)) \Big]$$

$$\leq \underset{(t,v) \leftarrow (TV)|_W}{\mathrm{E}} \Big[ \log \frac{1}{\Pr[W \wedge V = v | T = t]} \Big]$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t}) \Big] \qquad \text{(chain rule)}$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum D((X_j|\underbrace{W, V = v}_{w'})|_{T=t} \| (X_j|T=t)) \Big]$$

$$\leq \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \log \frac{1}{\Pr[W \wedge V = v|T=t]} \Big] \qquad \text{(Thm 8)}$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X'_j(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X'_j(T))$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t}) \Big] \qquad \text{(chain rule)}$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum D((X_j|\underbrace{W, \; V = v}_{w'})|_{T=t} \| (X_j|T=t)) \Big]$$

$$\leq \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \log \frac{1}{\Pr[W \wedge V = v | T = t]} \Big] \qquad \text{(Thm 8)}$$

$$\leq \log \mathop{E}_{(t,v) \leftarrow (TV)|_W} \frac{1}{\Pr[W \wedge V = v | T = t]}$$

## Proving **Thm 11**

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))$$

$$= \operatorname*{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum_{j=1}^{k} D(X_j|_{W, V=v, T=t} \| (X_j|_{T=t}) \Big] \qquad \text{(chain rule)}$$

$$= \operatorname*{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum D((X_j|\underbrace{W, V = v}_{w'})|_{T=t} \| (X_j|T = t)) \Big]$$

$$\leq \operatorname*{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \log \frac{1}{\Pr[W \wedge V = v|T = t]} \Big] \qquad \text{(Thm 8)}$$

$$\leq \log \operatorname*{E}_{(t,v) \leftarrow (TV)|_W} \frac{1}{\Pr[W \wedge V = v|T = t]} \qquad \text{(Jensen's inequality)}$$

## Proving Thm 11

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum_{j=1}^{k} D(X_j|_{W,V=v,T=t} \| (X_j|_{T=t}) \Big] \qquad \text{(chain rule)}$$

$$= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \sum D((X_j|_{\underbrace{W, V = v}_{w'}})|_{T=t} \| (X_j | T = t)) \Big]$$

$$\leq \mathop{E}_{(t,v) \leftarrow (TV)|_W} \Big[ \log \frac{1}{\Pr[W \wedge V = v | T = t]} \Big] \qquad \text{(Thm 8)}$$

$$\leq \log \mathop{E}_{(t,v) \leftarrow (TV)|_W} \frac{1}{\Pr[W \wedge V = v | T = t]} \qquad \text{(Jensen's inequality)}$$

$$= \log \sum_{(t,v) \in \text{Supp}((TV)|_W)} \frac{\Pr[T = t]}{\Pr[W]}$$

## Proving **Thm 11**

Let $X = (X_1, \ldots, X_k)$, $T$ and $V$ be rv's over $\mathcal{X}^k$, $\mathcal{T}$ and $\mathcal{V}$ respectively, such that $X_i$'s are iid conditioned on $T$. Let $W$ be an event and let $X_j'(t)$ be distributed according to the distribution of $X_j|_{T=t}$.

$$
\sum_{j=1}^{k} D((TVX_j)|_W \| (TV)|_W X_j'(T))
$$

$$
= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \left[ \sum_{j=1}^{k} D(X_j|_{W, V=v, T=t} \| (X_j|_{T=t}) \right] \qquad \text{(chain rule)}
$$

$$
= \mathop{E}_{(t,v) \leftarrow (TV)|_W} \left[ \sum D((X_j|\underbrace{_{W, V = v}}_{w'})|_{T=t} \| (X_j|T = t)) \right]
$$

$$
\leq \mathop{E}_{(t,v) \leftarrow (TV)|_W} \left[ \log \frac{1}{\Pr[W \wedge V = v | T = t]} \right] \qquad \text{(Thm 8)}
$$

$$
\leq \log \mathop{E}_{(t,v) \leftarrow (TV)|_W} \frac{1}{\Pr[W \wedge V = v | T = t]} \qquad \text{(Jensen's inequality)}
$$

$$
= \log \sum_{(t,v) \in \mathrm{Supp}((TV)|_W)} \frac{\Pr[T = t]}{\Pr[W]} \leq \log \frac{\|\mathrm{Supp}(V|_w)\|}{\Pr[W]}. \qquad \square
$$